

Mining generalized fuzzy association rules from web pages

Yi-Tsung Tang, Master Student

Hung-Pin Chiu, Assistant Professor

Department of Information Management, NAN HUA University

Department of Information Management, National Taitung University

eric_tang630@hotmail.com

hpchiu@nttu.edu.tw

Abstract

The discovery of fuzzy association rules is an important data-mining task for which many algorithms have been proposed. However, the efficiency of these algorithms needs to be improved to handle real-world large datasets. In this paper, we present an efficient method named cluster-based fuzzy association rule (CBFAR) to discover generalized fuzzy association rules from web pages. The CBFAR method is to create fuzzy cluster tables by scanning the browse information database (BIDB) once, and then clustering the browse records to the k -th cluster table, where the length of a record is k . The counts of the fuzzy regions are stored in the Fuzzy_Cluster Tables. This method requires less contrast to generate large itemsets. The CBFAR method is also discussed.

Keyword : Fuzzy data mining; association rules

摘要

模糊關聯法則的挖掘是資料挖掘(Data Mining)中一個重要的部分，也有許多的方法相繼被提出。然而，這些演算法對於處理實際資料上的效率仍然有改進的空間。本研究提出了一個有效率的方法（Cluster-Based Fuzzy Association Rule:CBFAR）來從許多網頁中找出模糊關聯法則，並改進挖掘的處理效率，此方法以分群表(cluster table)的觀念來儲存網頁瀏覽次數之模糊值，在大項目組的產生過程中，只需掃描瀏覽資料庫一次並去除許多不必要的資料比對時間，有效的減少處理時間，改進效率。

關鍵詞：模糊資料挖掘、關聯法則

1. Introduction

The discovery of fuzzy association rules is an important data-mining task. Association rules are used to discover the relationships, and potential associations, of items or attributes among huge data. These rules can be effective in uncovering unknown relationships, providing results that can be the basis of forecast and decision.

Deriving association rules from transaction database is most commonly seen in data mining. [2][4] It discovers relationships among items. In the past, Agrawal and Srikant proposed the Apriori association rule algorithm.[5] It can discover meaningful itemsets and construct association rules within large databases, but a large number of the candidate itemsets are generated from single itemsets. This method also needs to perform contrasts against all of the transactions, level by level, in the process of creating association rules. The database is repeatedly scanned to contrast each candidate itemset, that performance is dramatically affected.

After Agrawal et al. proposed the Apriori association rule, Tsay et al. have used cluster-based association rule (CBAR) approach.[8] This method used cluster-based table to reduce the number of database scans and requiring less contrast. Recently, the fuzzy set theory[3] has been used more and more frequently in intelligent systems. It's simplicity and similarity to human reasoning.[1] Hong et al. also proposed a fuzzy mining algorithm.[7] The items considered in their approach had hierarchical

relationships. However, items in real-world applications are usually organized in some hierarchies. Mining multiple-concept-level fuzzy rules may lead to discovery of more general and important knowledge from data.

In this paper, we present a new method called cluster-based fuzzy association rule (CBFAR), for efficient fuzzy association rules mining. We considered the hierarchical relationships to discover the generalized fuzzy association rules from the browse information database (BIDB) and used the cluster-based concept to reduce the number of database scans. When the customer clicking the web pages, then the click times stored in the browse information database (BIDB). This method not only needs only one database scans, but also requires less contrast.

2. CBFAR Mining Framework

The hierarchical relationships and cluster-based concepts are used to discover generalized fuzzy association rules from browse information database (BIDB). We propose a CBFAR mining framework for discovering generalized fuzzy association rules. The proposed framework is shown in Fig. 1.

We proposed mining framework maintains fuzzy association rules, and uses the hierarchical relationships and cluster-based fuzzy table to derive the fuzzy association rules. Previous studies on data mining focused on finding association rules on the single-concept level. However, relevant web page taxonomies are usually predefined in the networks structure and can be represented using hierarchical trees.[6] Terminal nodes on the trees represent actual web pages appearing in networks structure; internal nodes represent main or sub-main web pages formed by lower-level nodes. A simple example is given in Fig. 2.

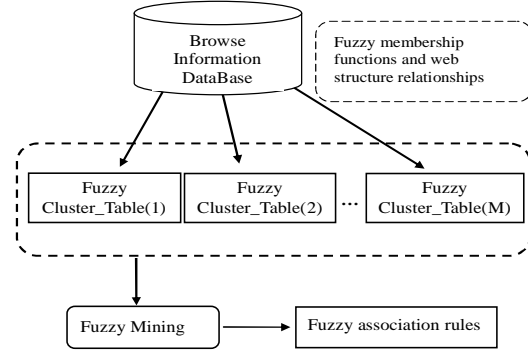


Figure1: CBFAR Mining Framework

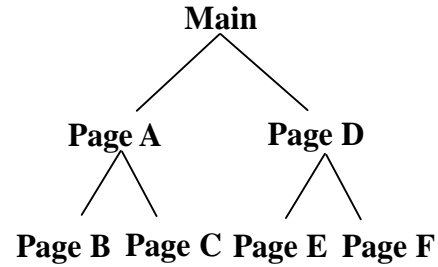


Figure2: An example of taxonomic structures

In this example, the main page falls into two sub-pages: page A and page D. Page A can be further classified into page B and page C. Similarly, assume page D are divided into page E and page F. The web pages (A, B, C, D and E) can appear in browse information records. The CBFAR mining method is divided into four phases.

In the first phase, all of web pages in each given browsed records are added according to the predefined taxonomy.

In the second phase, transform the quantitative value v_{ij} of each browsed data D_i ($i=1$ to n), for each expanded item name I_j appearing into a fuzzy set f_{ij} . The f_{ij} are represented as $(f_{ij1}/R_{j1} + f_{ij2}/R_{j2} + \dots + f_{ijh}/R_{jh})$ using the given membership functions, where h is the number of fuzzy regions for I_j . R_{jl} is the l th fuzzy region of I_j , $1 \leq l \leq h$, and f_{ijl} is v_{ij} 's fuzzy membership value in region R_{jl} . Calculate the value of each fuzzy region R_{jl} in the browsed data. $(count_{jl} = \sum_{i=1}^n f_{ijl})$

In the third phase, creates M cluster tables. Scan the browse information database once and cluster the

browsed data. If the length of browsed data is k , the browsed record and the fuzzy region value of items in this browsed record will be stored in the table, named Fuzzy_Cluster Table (k), $1 \leq k \leq M$, where M is the length of the longest browsed record in database.

In the fourth phase, the set of candidate itemsets C_n is generated. When the length of candidate itemset is k , the support is calculated with reference to the Fuzzy_Cluster Table(k). If the fuzzy region value of C_n is greater than or equal to the predefined minimum support value α , the candidate itemsets becomes the large itemsets, put C_n in the large itemsets L_n . Otherwise, it is contrasted with the Fuzzy_Cluster Table($k+1$). The large itemsets is $L_n = \{ \max-R_j \mid \max-count_j \geq \alpha, 1 \leq j \leq m \}$. Until the large itemsets L_n is null, this process terminates when the calculated support is greater than or equal to the predefined minimum support or the the end of the Fuzzy_Cluster Table(M) has been reached. Finally, use the predefined minimum confidence value to discover fuzzy association rules. If the candidate fuzzy association rule is larger than or equal to the predefined confidence value, put it in the rule base.

3. An Example

In this section, an example is given to illustrate the proposed mining method. This is a simple example to show how the proposed method can be used to discover fuzzy association rules from browsed data. There are six browsed records and five items (web pages) in a browse information database: A, B, C, D, E and F. An example browse information database is shown in Table 1. The taxonomy tree is shown in Fig. 3. All of items (web pages) appearing in the browse information database (BIDB) according to the predefined taxonomy tree.

Table 1. Six browsed records in this example

BID	Items (Web Pages, Click times)
B1	(A,3) (B,4) (C,2) (D,3) (E,4) (F,2)
B2	(A,3) (B,7) (C,7) (D,3) (E,7) (F,7)
B3	(A,4) (B,2) (C,5) (E,6) (F,5)
B4	(B,9) (C,10) (D,9) (E,10)
B5	(B,3) (F,3)
B6	(B,8)(D,4) (E,8) (F,4)

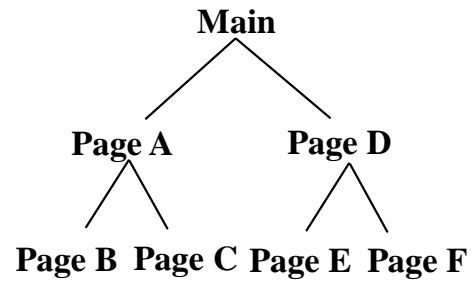


Figure3: Taxonomy tree in this example

In this example, assume that the fuzzy membership functions are the same for all the items and are as shown in Fig. 4. The fuzzy membership function is represented by three fuzzy regions: *Low(L)*, *Middle(M)* and *High(H)*, and three fuzzy membership values are produced for each item according to the predefined membership function.

The length of the longest browsed record in this database is six, and creates six fuzzy_cluster tables as shown in Table 2. The fuzzy region value of items in this browsed record will be stored in the Fuzzy_Cluster Tables.

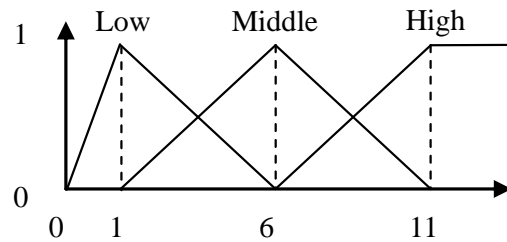


Figure4: The membership function in this example

Table 2: Fuzzy_Cluster Tables

BID	A	B	C	D	E	F
Fuzzy_Cluster Table(1)						
NULL						
Fuzzy_Cluster Table(2)						
B5	0	L,0.6 M,0.4	0	0	0	L,0.6 M,0.4
Fuzzy_Cluster Table(3)						
NULL						
Fuzzy_Cluster Table(4)						
B4	0	M,0.4 H,0.6	M,0.2 H,0.8	M,0.4 H,0.6	M,0.2 H,0.8	0
B6	0	M,0.6 H,0.4	0	L,0.4 M,0.6	M,0.6 H,0.4	L,0.4 M,0.6
Fuzzy_Cluster Table(5)						
B3	L,0.4 M,0.6	L,0.8 M,0.2	L,0.2 M,0.8	0	M,1.0	L,0.2 M,0.8
Fuzzy_Cluster Table(6)						
B1	L,0.6 M,0.4	L,0.4 M,0.6	L,0.8 M,0.2	L,0.6 M,0.4	L,0.4 M,0.6	L,0.8 M,0.2
B2	L,0.6 M,0.4	M,0.8 H,0.2	M,0.8 H,0.2	L,0.6 M,0.4	M,0.8 H,0.2	M,0.8 H,0.2

Assume the minimum support value is 2.4. We can discover the Large-1 itemsets (L_1) which is large than or equal to the predefined minimum support value according to the fuzzy_cluster tables. The itemsets of L_1 are {B.Middle = 3.0}, {E.Middle = 3.2}, {F.Middle = 2.8}.

Generate the large 2-itemsets L_2 . Combining the items of L_1 in order to generate candidate 2-itemsets C_2 . The procedure is similar to the candidate generation of Apriori algorithm[5]. The itemsets of C_2 are {B.Middle,E.Middle}, {B.Middle,F.Middle}, {E.Middle,F.Middle}. In order to generate L_2 , it is necessary to compute the fuzzy region values of each candidate itemset in the Fuzzy_Cluster Table(2). If the value is larger than or equal to the predefined minimum support value, put C_2 in the L_2 . Otherwise, compute the fuzzy region values in the next cluster table (Fuzzy_Cluster Table(3)). The other large itemsets L_n are in the similar way.

Therefore, the large itemsets in this example are {B.Middle},{E.Middle},{F.Middle},{B.Middle,E.Middle},{E.Middle,F.Middle}. Then, we can transform

each large itemsets into a fuzzy association rule. In the electronic commerce (EC) environment, we can use the association rules to fascinate the customer. Then the customer relationship management (CRM) can make a better profit.

4. Conclusions

In this paper, we have proposed a generalized fuzzy association rules mining framework for extracting fuzzy association rules from browse information database (BIDB). In the electronic commerce environment, we can use the association rules to fascinate the customer. Then the customer relationship management (CRM) can make a better profit.

The cluster-based fuzzy association rule (CBFAR) method creates Fuzzy_Cluster Tables to discover the large itemsets. Contrasts are performed only against the partial Fuzzy_Cluster Tables that were created in advance. It only requires a single scan of the browse information database, and contrasts with the partial Fuzzy_Cluster Tables. This method not only needs only one database scans, but also requires less contrast.

In the future, we will continuously for the huge database, and discussing with the performance of CBFAR method.

5. References

- [1] A. Kandel, Fuzzy Expert Systems, CRC Press, Boca Raton, FL, 1992, pp. 8-19.
- [2] J. Han, Y. Fu, Discovery of multiple-level association rules from large database, The Internet. Conf. on Very Large Databases, 1995.
- [3] L.A. Zadeh, Fuzzy sets, Inform. and Control 8(3), 1965 pp. 338-353.
- [4] R.Agrawal, T. Imielinski, A. Swami, Mining

- association rules between sets of items in large database, the 1993 ACM SIGMOD Conf., Washington, DC, USA, 1993.
- [5] R. Agrawal, R. Srikant, Fast algorithm for mining association rules in large databasaes, Proceedings of 1994 International Conference on VLDB, 1994 pp. 487-499.
- [6] R. Srikant, R. Agrawal, Mining generalized association rules, The Internat. Conf. on Very Large Databases, 1995.
- [7] Tzung-Pei Hong, Kuei-Ying Lin, Shyue-Liang Wang, Fuzzy data mining for interesting generalized association rules, Fuzzy Sets and Systems, 2003 pp. 255-269.
- [8] Yuh-Jiuan Tsay, Jiunn-Yann Chiang, CBAR: an efficient method for mining association rules, Knowledge-Based Systems, 2005 pp. 99-105.